

Introduction to Machine Learning

Introduction to Linear Classifiers

Andres Mendez-Vazquez

February 11, 2019

Outline

1 Introduction

- Introduction
- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector

2 Developing a Solution

- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation
- Geometric Interpretation
- Issues with Least Squares!!!
 - Problem with Outliers
 - Problem with High Number of Dimensions

Outline

1 Introduction

● Introduction

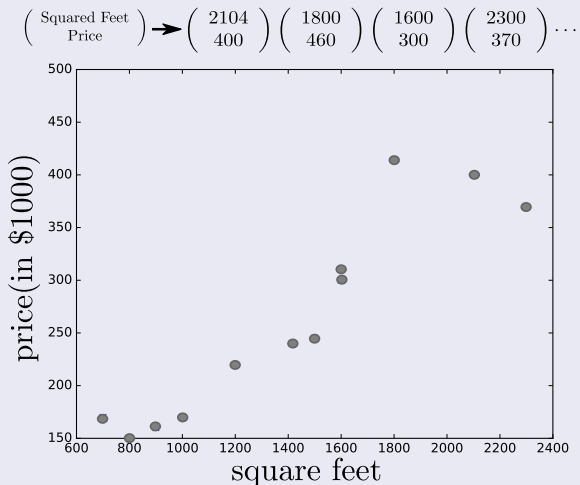
- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector

2 Developing a Solution

- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation
- Geometric Interpretation
- Issues with Least Squares!!!
 - Problem with Outliers
 - Problem with High Number of Dimensions

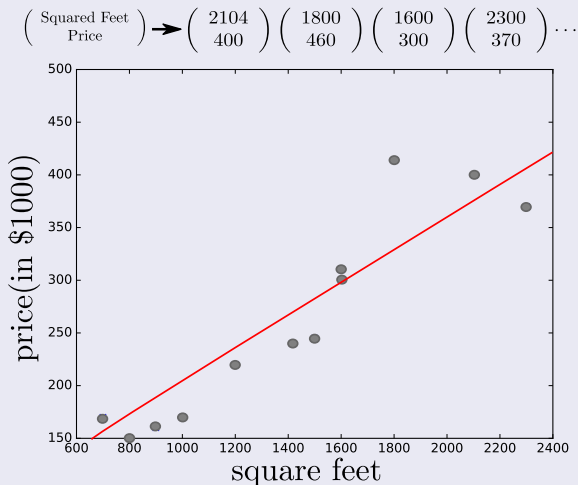
Many Times

We have this kind of data sets (House Prices)



Thus

We can adjust a line/hyperplane to be able to forecast prices



Thus, Our Objective

To find such hyperplane

To do forecasting on the prices of a house given its surface!!!

Here, where "Learning" Machine Learning style comes around

Basically, the process defined in Machine Learning!!!

Thus, Our Objective

To find such hyperplane

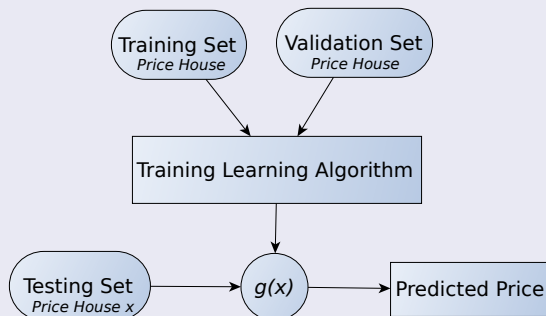
To do forecasting on the prices of a house given its surface!!!

Here, where “Learning” Machine Learning style comes around

Basically, the process defined in Machine Learning!!!

Then, in Supervised Training

We have the following process



Outline

1 Introduction

- Introduction
- **The Simplest Functions**
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector

2 Developing a Solution

- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation
- Geometric Interpretation
- Issues with Least Squares!!!
 - Problem with Outliers
 - Problem with High Number of Dimensions

What is it?

First than anything, we have a parametric model!!!

Here, we have an hyperplane as a model:

$$g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 \quad (1)$$

Note: $\mathbf{w}^T \mathbf{x}$ is also know as dot product

In the case of \mathbb{R}^2

We have:

$$g(\mathbf{x}) = (w_1, w_2) \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} + w_0 = w_1 x_1 + w_2 x_2 + w_0 \quad (2)$$

What is it?

First than anything, we have a parametric model!!!

Here, we have an hyperplane as a model:

$$g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 \quad (1)$$

Note: $\mathbf{w}^T \mathbf{x}$ is also know as dot product

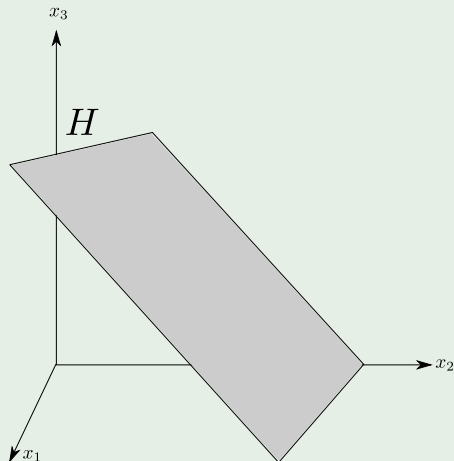
In the case of \mathbb{R}^2

We have:

$$g(\mathbf{x}) = (w_1, w_2) \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} + w_0 = w_1 x_1 + w_2 x_2 + w_0 \quad (2)$$

Example

Hyperplane in \mathbb{R}^3



Outline

1 Introduction

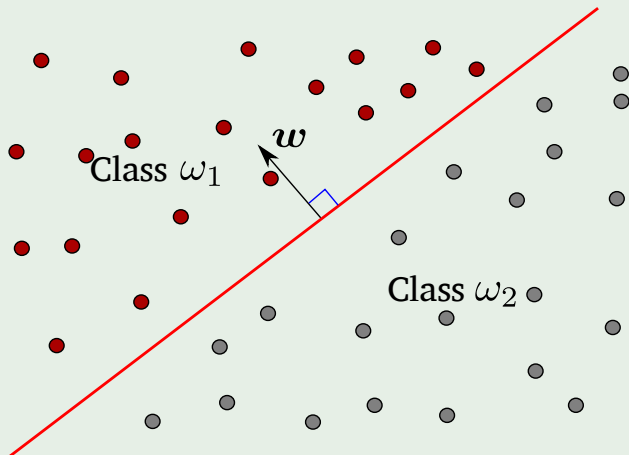
- Introduction
- The Simplest Functions
- **Splitting the Space**
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector

2 Developing a Solution

- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation
- Geometric Interpretation
- Issues with Least Squares!!!
 - Problem with Outliers
 - Problem with High Number of Dimensions

Splitting The Space \mathbb{R}^2

Using a simple straight line (Hyperplane)



Splitting the Space?

For example, assume the following vector w and constant w_0

$$w = (-1, 2)^T \text{ and } w_0 = 0$$

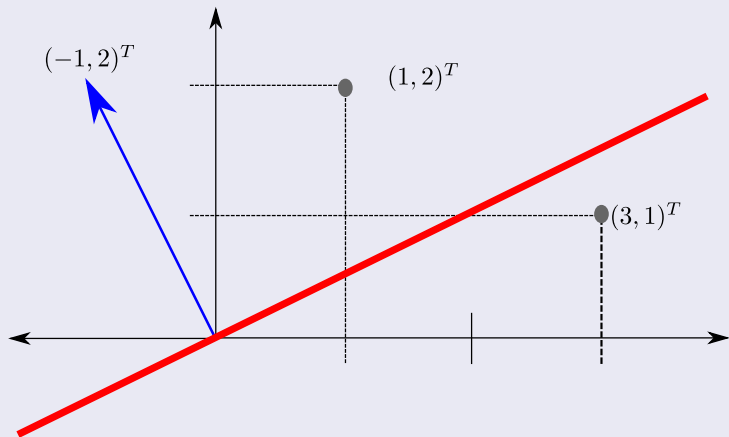
Hyperplane

Splitting the Space?

For example, assume the following vector w and constant w_0

$$w = (-1, 2)^T \text{ and } w_0 = 0$$

Hyperplane



Then, we have

The following results

$$g\left(\begin{pmatrix} 1 \\ 2 \end{pmatrix}\right) = (-1, 2) \begin{pmatrix} 1 \\ 2 \end{pmatrix} = -1 \times 1 + 2 \times 2 = 3$$

$$g\left(\begin{pmatrix} 3 \\ 1 \end{pmatrix}\right) = (-1, 2) \begin{pmatrix} 3 \\ 1 \end{pmatrix} = -1 \times 3 + 2 \times 1 = -1$$

YES!!! We have a positive side and a negative side!!!

Outline

1 Introduction

- Introduction
- The Simplest Functions
- Splitting the Space
- **Defining the Decision Surface**
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector

2 Developing a Solution

- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation
- Geometric Interpretation
- Issues with Least Squares!!!
 - Problem with Outliers
 - Problem with High Number of Dimensions

The Decision Surface

The equation $g(x) = 0$ defines a decision surface

Separating the elements in classes, ω_1 and ω_2 .

When $g(x)$ is linear the decision surface is an hyperplane

Now assume x_1 and x_2 are both on the decision surface

$$w^T x_1 + w_0 = 0$$

$$w^T x_2 + w_0 = 0$$

Thus

$$w^T x_1 + w_0 = w^T x_2 + w_0 \quad (3)$$

The Decision Surface

The equation $g(x) = 0$ defines a decision surface

Separating the elements in classes, ω_1 and ω_2 .

When $g(x)$ is linear the decision surface is an hyperplane

Now assume x_1 and x_2 are both on the decision surface

$$w^T x_1 + w_0 = 0$$

$$w^T x_2 + w_0 = 0$$

Thus

$$w^T x_1 + w_0 = w^T x_2 + w_0 \quad (3)$$

The Decision Surface

The equation $g(x) = 0$ defines a decision surface

Separating the elements in classes, ω_1 and ω_2 .

When $g(x)$ is linear the decision surface is an hyperplane

Now assume x_1 and x_2 are both on the decision surface

$$w^T x_1 + w_0 = 0$$

$$w^T x_2 + w_0 = 0$$

Thus

$$w^T x_1 + w_0 = w^T x_2 + w_0 \quad (3)$$

Defining a Decision Surface

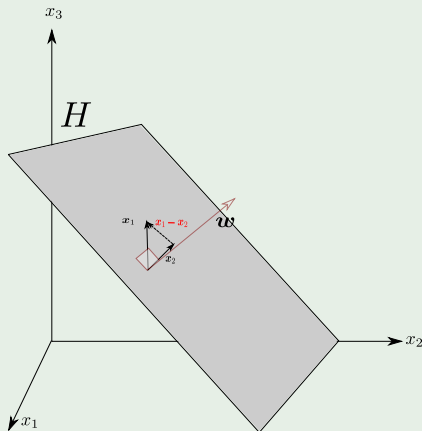
Then

$$\mathbf{w}^T (\mathbf{x}_1 - \mathbf{x}_2) = 0 \quad (4)$$

Therefore

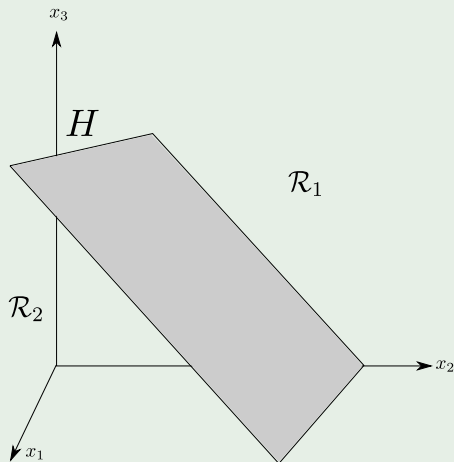
$x_1 - x_2$ lives in the hyperplane i.e. it is perpendicular to w^T

- Remark: any vector in the hyperplane is a linear combination of elements in a basis
- **Therefore any vector in the plane is perpendicular to w^T**



Therefore

The space is split in two regions (Example in \mathbb{R}^3) by the hyperplane H



Outline

1 Introduction

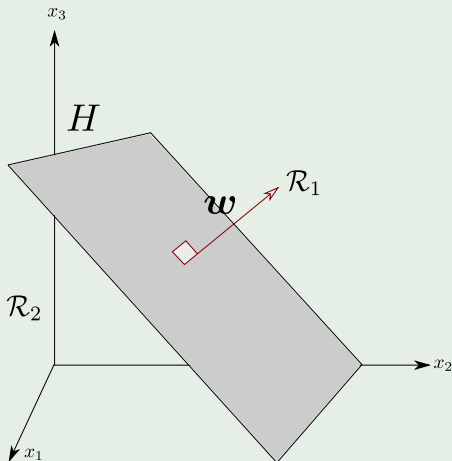
- Introduction
- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- **Properties of the Hyperplane $w^T x + w_0$**
- Augmenting the Vector

2 Developing a Solution

- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation
- Geometric Interpretation
- Issues with Least Squares!!!
 - Problem with Outliers
 - Problem with High Number of Dimensions

Some Properties of the Hyperplane

Given that $g(\mathbf{x}) > 0$ if $\mathbf{x} \in \mathcal{R}_1$



It is more

We can say the following

- Any $x \in \mathcal{R}_1$ is on the positive side of H .
- Any $x \in \mathcal{R}_2$ is on the negative side of H .

It is more

We can say the following

- Any $x \in \mathcal{R}_1$ is on the positive side of H .
- Any $x \in \mathcal{R}_2$ is on the negative side of H .

In addition, you can give us a way to obtain the distance from x to the hyperplane H .

First, we express any x as follows

$$x = x_p + r \frac{w}{\|w\|}$$

It is more

We can say the following

- Any $x \in \mathcal{R}_1$ is on the positive side of H .
- Any $x \in \mathcal{R}_2$ is on the negative side of H .

In addition, $g(x)$ can give us a way to obtain the distance from x to the hyperplane H

First, we express any x as follows

$$x = x_p + r \frac{w}{\|w\|}$$

Where

- x_p is the normal projection of x onto H .
- r is the desired distance
 - ▶ Positive, if x is in the positive side
 - ▶ Negative, if x is in the negative side

It is more

We can say the following

- Any $x \in \mathcal{R}_1$ is on the positive side of H .
- Any $x \in \mathcal{R}_2$ is on the negative side of H .

In addition, $g(x)$ can give us a way to obtain the distance from x to the hyperplane H

First, we express any x as follows

$$x = x_p + r \frac{w}{\|w\|}$$

Where

- x_p is the normal projection of x onto H .
- r is the desired distance
 - Positive, if x is in the positive side
 - Negative, if x is in the negative side

It is more

We can say the following

- Any $x \in \mathcal{R}_1$ is on the positive side of H .
- Any $x \in \mathcal{R}_2$ is on the negative side of H .

In addition, $g(x)$ can give us a way to obtain the distance from x to the hyperplane H

First, we express any x as follows

$$x = x_p + r \frac{w}{\|w\|}$$

Where

- x_p is the normal projection of x onto H .
- r is the desired distance

► Positive, if x is in the positive side
► Negative, if x is in the negative side

It is more

We can say the following

- Any $x \in \mathcal{R}_1$ is on the positive side of H .
- Any $x \in \mathcal{R}_2$ is on the negative side of H .

In addition, $g(x)$ can give us a way to obtain the distance from x to the hyperplane H

First, we express any x as follows

$$x = x_p + r \frac{w}{\|w\|}$$

Where

- x_p is the normal projection of x onto H .
- r is the desired distance
 - ▶ Positive, if x is in the positive side
 - ▶ Negative, if x is in the negative side

It is more

We can say the following

- Any $x \in \mathcal{R}_1$ is on the positive side of H .
- Any $x \in \mathcal{R}_2$ is on the negative side of H .

In addition, $g(x)$ can give us a way to obtain the distance from x to the hyperplane H

First, we express any x as follows

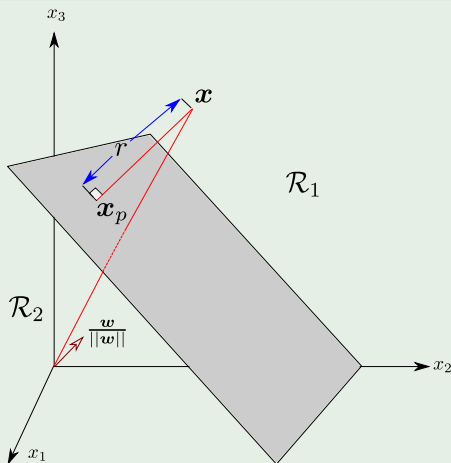
$$x = x_p + r \frac{w}{\|w\|}$$

Where

- x_p is the normal projection of x onto H .
- r is the desired distance
 - ▶ Positive, if x is in the positive side
 - ▶ Negative, if x is in the negative side

We have something like this

We have then



Now

Since $g(x_p) = 0$

We have that

$$\begin{aligned} g(x) &= g\left(x_p + r \frac{w}{\|w\|}\right) \\ &= w^T \left(x_p + r \frac{w}{\|w\|}\right) + w_0 \\ &= w^T x_p + w_0 + r \frac{w^T w}{\|w\|} \\ &= g(x_p) + r \frac{\|w\|^2}{\|w\|} \\ &= r \|w\| \end{aligned}$$

Now

Since $g(\mathbf{x}_p) = 0$

We have that

$$\begin{aligned} g(\mathbf{x}) &= g\left(\mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|}\right) \\ &= \mathbf{w}^T \left(\mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|}\right) + w_0 \\ &= \mathbf{w}^T \mathbf{x}_p + w_0 + r \frac{\mathbf{w}^T \mathbf{w}}{\|\mathbf{w}\|} \\ &= g(\mathbf{x}_p) + r \frac{\|\mathbf{w}\|^2}{\|\mathbf{w}\|} \\ &= r \|\mathbf{w}\| \end{aligned}$$

Then, we have

$$r = \frac{g(\mathbf{x})}{\|\mathbf{w}\|} \quad (5)$$

Now

Since $g(\mathbf{x}_p) = 0$

We have that

$$\begin{aligned} g(\mathbf{x}) &= g\left(\mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|}\right) \\ &= \mathbf{w}^T \left(\mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|}\right) + w_0 \\ &= \mathbf{w}^T \mathbf{x}_p + w_0 + r \frac{\mathbf{w}^T \mathbf{w}}{\|\mathbf{w}\|} \\ &= g(\mathbf{x}_p) + r \frac{\|\mathbf{w}\|^2}{\|\mathbf{w}\|} \\ &= r \|\mathbf{w}\| \end{aligned}$$

Then, we have

$$r = \frac{g(\mathbf{x})}{\|\mathbf{w}\|} \quad (5)$$

Now

Since $g(\mathbf{x}_p) = 0$

We have that

$$\begin{aligned}g(\mathbf{x}) &= g\left(\mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|}\right) \\&= \mathbf{w}^T \left(\mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|}\right) + w_0 \\&= \mathbf{w}^T \mathbf{x}_p + w_0 + r \frac{\mathbf{w}^T \mathbf{w}}{\|\mathbf{w}\|} \\&= g(\mathbf{x}_p) + r \frac{\|\mathbf{w}\|^2}{\|\mathbf{w}\|} \\&= r \|\mathbf{w}\|\end{aligned}$$

Then, we have

$$r = \frac{g(\mathbf{x})}{\|\mathbf{w}\|} \quad (5)$$

Now

Since $g(\mathbf{x}_p) = 0$

We have that

$$\begin{aligned} g(\mathbf{x}) &= g\left(\mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|}\right) \\ &= \mathbf{w}^T \left(\mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|}\right) + w_0 \\ &= \mathbf{w}^T \mathbf{x}_p + w_0 + r \frac{\mathbf{w}^T \mathbf{w}}{\|\mathbf{w}\|} \\ &= g(\mathbf{x}_p) + r \frac{\|\mathbf{w}\|^2}{\|\mathbf{w}\|} \\ &= r \|\mathbf{w}\| \end{aligned}$$

Then, we have

$$r = \frac{g(\mathbf{x})}{\|\mathbf{w}\|}$$

(5)

Now

Since $g(\mathbf{x}_p) = 0$

We have that

$$\begin{aligned}g(\mathbf{x}) &= g\left(\mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|}\right) \\&= \mathbf{w}^T \left(\mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|}\right) + w_0 \\&= \mathbf{w}^T \mathbf{x}_p + w_0 + r \frac{\mathbf{w}^T \mathbf{w}}{\|\mathbf{w}\|} \\&= g(\mathbf{x}_p) + r \frac{\|\mathbf{w}\|^2}{\|\mathbf{w}\|} \\&= r \|\mathbf{w}\|\end{aligned}$$

Then, we have

$$r = \frac{g(\mathbf{x})}{\|\mathbf{w}\|} \quad (5)$$

In particular

The distance from the origin to H

$$r = \frac{g(\mathbf{0})}{\|\mathbf{w}\|} = \frac{\mathbf{w}^T(\mathbf{0}) + w_0}{\|\mathbf{w}\|} = \frac{w_0}{\|\mathbf{w}\|} \quad (6)$$

In particular

The distance from the origin to H

$$r = \frac{g(\mathbf{0})}{\|\mathbf{w}\|} = \frac{\mathbf{w}^T(\mathbf{0}) + w_0}{\|\mathbf{w}\|} = \frac{w_0}{\|\mathbf{w}\|} \quad (6)$$

Remarks

- If $w_0 > 0$, the origin is on the positive side of H .
- If $w_0 < 0$, the origin is on the negative side of H .
- If $w_0 = 0$, the hyperplane has the homogeneous form $\mathbf{w}^T \mathbf{x}$ and hyperplane passes through the origin.

In particular

The distance from the origin to H

$$r = \frac{g(\mathbf{0})}{\|\mathbf{w}\|} = \frac{\mathbf{w}^T(\mathbf{0}) + w_0}{\|\mathbf{w}\|} = \frac{w_0}{\|\mathbf{w}\|} \quad (6)$$

Remarks

- If $w_0 > 0$, the origin is on the positive side of H .
- If $w_0 < 0$, the origin is on the negative side of H .
- If $w_0 = 0$, the hyperplane has the homogeneous form $\mathbf{w}^T \mathbf{x}$ and hyperplane passes through the origin.

In particular

The distance from the origin to H

$$r = \frac{g(\mathbf{0})}{\|\mathbf{w}\|} = \frac{\mathbf{w}^T(\mathbf{0}) + w_0}{\|\mathbf{w}\|} = \frac{w_0}{\|\mathbf{w}\|} \quad (6)$$

Remarks

- If $w_0 > 0$, the origin is on the positive side of H .
- If $w_0 < 0$, the origin is on the negative side of H .
- If $w_0 = 0$, the hyperplane has the homogeneous form $\mathbf{w}^T \mathbf{x}$ and hyperplane passes through the origin.

Outline

1 Introduction

- Introduction
- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- **Augmenting the Vector**

2 Developing a Solution

- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation
- Geometric Interpretation
- Issues with Least Squares!!!
 - Problem with Outliers
 - Problem with High Number of Dimensions

We want to solve the independence of w_0

We would like w_0 as part of the dot product by making $x_0 = 1$

$$g(\mathbf{x}) = w_0 \times 1 + \sum_{i=1}^d w_i x_i = w_0 \times x_0 + \sum_{i=1}^d w_i x_i = \sum_{i=0}^d w_i x_i \quad (7)$$

We want to solve the independence of w_0

We would like w_0 as part of the dot product by making $x_0 = 1$

$$g(\mathbf{x}) = w_0 \times 1 + \sum_{i=1}^d w_i x_i = w_0 \times x_0 + \sum_{i=1}^d w_i x_i = \sum_{i=0}^d w_i x_i \quad (7)$$

By making

$$\mathbf{x}_{aug} = \begin{pmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{pmatrix} = \begin{pmatrix} 1 \\ \mathbf{x} \end{pmatrix}$$

We want to solve the independence of w_0

We would like w_0 as part of the dot product by making $x_0 = 1$

$$g(\mathbf{x}) = w_0 \times 1 + \sum_{i=1}^d w_i x_i = w_0 \times x_0 + \sum_{i=1}^d w_i x_i = \sum_{i=0}^d w_i x_i \quad (7)$$

By making

$$\mathbf{x}_{aug} = \begin{pmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{pmatrix} = \begin{pmatrix} 1 \\ \mathbf{x} \end{pmatrix}$$

Where

\mathbf{x}_{aug} is called an augmented feature vector.

We want to solve the independence of w_0

We would like w_0 as part of the dot product by making $x_0 = 1$

$$g(\mathbf{x}) = w_0 \times 1 + \sum_{i=1}^d w_i x_i = w_0 \times x_0 + \sum_{i=1}^d w_i x_i = \sum_{i=0}^d w_i x_i \quad (7)$$

By making

$$\mathbf{x}_{aug} = \begin{pmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{pmatrix} = \begin{pmatrix} 1 \\ \mathbf{x} \end{pmatrix}$$

Where

\mathbf{x}_{aug} is called an augmented feature vector.

We want to solve the independence of w_0

We would like w_0 as part of the dot product by making $x_0 = 1$

$$g(\mathbf{x}) = w_0 \times 1 + \sum_{i=1}^d w_i x_i = w_0 \times x_0 + \sum_{i=1}^d w_i x_i = \sum_{i=0}^d w_i x_i \quad (7)$$

By making

$$\mathbf{x}_{aug} = \begin{pmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{pmatrix} = \begin{pmatrix} 1 \\ \mathbf{x} \end{pmatrix}$$

Where

\mathbf{x}_{aug} is called an augmented feature vector.

In a similar way

We have the augmented weight vector

$$\mathbf{w}_{aug} = \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{pmatrix} = \begin{pmatrix} w_0 \\ \mathbf{w} \end{pmatrix}$$

In a similar way

We have the augmented weight vector

$$\mathbf{w}_{aug} = \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{pmatrix} = \begin{pmatrix} w_0 \\ \mathbf{w} \end{pmatrix}$$

Remarks

- The addition of a constant component to \mathbf{x} preserves all the distance relationship between samples.
- The resulting \mathbf{x}_{aug} vectors, all lie in a d -dimensional subspace which is the \mathbf{x} -space itself.

In a similar way

We have the augmented weight vector

$$\mathbf{w}_{aug} = \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{pmatrix} = \begin{pmatrix} w_0 \\ \mathbf{w} \end{pmatrix}$$

Remarks

- The addition of a constant component to \mathbf{x} preserves all the distance relationship between samples.
- The resulting \mathbf{x}_{aug} vectors, all lie in a d -dimensional subspace which is the \mathbf{x} -space itself.

Outline

1 Introduction

- Introduction
- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector

2 Developing a Solution

- **Least Squared Error Procedure**
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation
- Geometric Interpretation
- Issues with Least Squares!!!
 - Problem with Outliers
 - Problem with High Number of Dimensions

Outline

1 Introduction

- Introduction
- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector

2 Developing a Solution

- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation
- Geometric Interpretation
- Issues with Least Squares!!!
 - Problem with Outliers
 - Problem with High Number of Dimensions

Initial Supposition

Suppose, we have

n samples $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ some labeled ω_1 and some labeled ω_2 .

Initial Supposition

Suppose, we have

n samples $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ some labeled ω_1 and some labeled ω_2 .

We want a vector weight \mathbf{w} such that

- $\mathbf{w}^T \mathbf{x}_i > 0$, if $\mathbf{x}_i \in \omega_1$.

- $\mathbf{w}^T \mathbf{x}_i < 0$, if $\mathbf{x}_i \in \omega_2$.

The name of this weight vector

It is called a separating vector or solution vector.

Initial Supposition

Suppose, we have

n samples $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ some labeled ω_1 and some labeled ω_2 .

We want a vector weight \mathbf{w} such that

- $\mathbf{w}^T \mathbf{x}_i > 0$, if $\mathbf{x}_i \in \omega_1$.
- $\mathbf{w}^T \mathbf{x}_i < 0$, if $\mathbf{x}_i \in \omega_2$.

The name of this weight vector

It is called a separating vector or solution vector.

Initial Supposition

Suppose, we have

n samples $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ some labeled ω_1 and some labeled ω_2 .

We want a vector weight \mathbf{w} such that

- $\mathbf{w}^T \mathbf{x}_i > 0$, if $\mathbf{x}_i \in \omega_1$.
- $\mathbf{w}^T \mathbf{x}_i < 0$, if $\mathbf{x}_i \in \omega_2$.

The name of this weight vector

It is called a separating vector or solution vector.

Now, assume the following

Imagine that your problem has two classes ω_1 and ω_2 in \mathbb{R}^2

① They are linearly separable!!!

② You require to label them.

Now, assume the following

Imagine that your problem has two classes ω_1 and ω_2 in \mathbb{R}^2

- ① They are linearly separable!!!
- ② You require to label them.

We have a problem!!!

Which is the problem?

Now, assume the following

Imagine that your problem has two classes ω_1 and ω_2 in \mathbb{R}^2

- ① They are linearly separable!!!
- ② You require to label them.

We have a problem!!!

Which is the problem?

We do not know the hyperplane!!!

Thus, what distance each point has to the hyperplane?

Now, assume the following

Imagine that your problem has two classes ω_1 and ω_2 in \mathbb{R}^2

- ① They are linearly separable!!!
- ② You require to label them.

We have a problem!!!

Which is the problem?

We do not know the hyperplane!!!

Thus, what distance each point has to the hyperplane?

A Simple Solution For Our Quandary

Label the Classes

- $\omega_1 \implies +1$
- $\omega_2 \implies -1$

A Simple Solution For Our Quandary

Label the Classes

- $\omega_1 \implies +1$
- $\omega_2 \implies -1$

We produce the following labels

- 1 if $\mathbf{x} \in \omega_1$ then $y_{ideal} = g_{ideal}(\mathbf{x}) = +1$.
- 2 if $\mathbf{x} \in \omega_2$ then $y_{ideal} = g_{ideal}(\mathbf{x}) = -1$.

Remark: We have a problem with this labels!!!

A Simple Solution For Our Quandary

Label the Classes

- $\omega_1 \implies +1$
- $\omega_2 \implies -1$

We produce the following labels

- 1 if $\mathbf{x} \in \omega_1$ then $y_{ideal} = g_{ideal}(\mathbf{x}) = +1$.
- 2 if $\mathbf{x} \in \omega_2$ then $y_{ideal} = g_{ideal}(\mathbf{x}) = -1$.

Remark: We have a problem with this labels!!!

A Simple Solution For Our Quandary

Label the Classes

- $\omega_1 \implies +1$
- $\omega_2 \implies -1$

We produce the following labels

- 1 if $\mathbf{x} \in \omega_1$ then $y_{ideal} = g_{ideal}(\mathbf{x}) = +1$.
- 2 if $\mathbf{x} \in \omega_2$ then $y_{ideal} = g_{ideal}(\mathbf{x}) = -1$.

Remark: We have a problem with this labels!!!

Outline

1 Introduction

- Introduction
- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector

2 Developing a Solution

- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- **The Error Idea**
- The Final Error Equation
- Geometric Interpretation
- Issues with Least Squares!!!
 - Problem with Outliers
 - Problem with High Number of Dimensions

Now, What?

Assume true function f is given by

$$y_{noise} = g_{noise}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 + e \quad (8)$$

Where the e

It has a $e \sim N(\mu, \sigma^2)$

Thus, we can do the following

$$y_{noise} = g_{noise}(\mathbf{x}) = g_{ideal}(\mathbf{x}) + e \quad (9)$$

Now, What?

Assume true function f is given by

$$y_{noise} = g_{noise}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 + e \quad (8)$$

Where the e

It has a $e \sim N(\mu, \sigma^2)$

Thus, we can do the following

$$y_{noise} = g_{noise}(\mathbf{x}) = g_{ideal}(\mathbf{x}) + e \quad (9)$$

Now, What?

Assume true function f is given by

$$y_{noise} = g_{noise}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 + e \quad (8)$$

Where the e

It has a $e \sim N(\mu, \sigma^2)$

Thus, we can do the following

$$y_{noise} = g_{noise}(\mathbf{x}) = g_{ideal}(\mathbf{x}) + e \quad (9)$$

Thus, we have

What to do?

$$e = y_{noise} - g_{ideal}(\mathbf{x}) \quad (10)$$

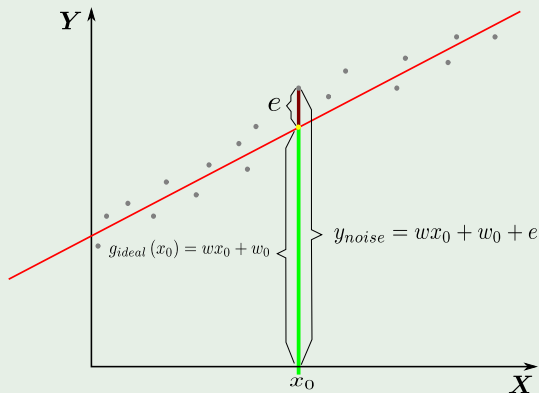
Graphically

Thus, we have

What to do?

$$e = y_{\text{noise}} - g_{\text{ideal}}(\mathbf{x}) \quad (10)$$

Graphically



Then, we have

A TRICK... Quite a good one!!! Instead of using y_{noise}

$$e = y_{noise} - g_{ideal}(\mathbf{x}) \quad (11)$$

We use y_{ideal}

$$e = y_{ideal} - g_{ideal}(\mathbf{x}) \quad (12)$$

We will see

How the geometry will solve the problem with using these labels.

Then, we have

A TRICK... Quite a good one!!! Instead of using y_{noise}

$$e = y_{noise} - g_{ideal}(\mathbf{x}) \quad (11)$$

We use y_{ideal}

$$e = y_{ideal} - g_{ideal}(\mathbf{x}) \quad (12)$$

We will see

How the geometry will solve the problem with using these labels.

Then, we have

A TRICK... Quite a good one!!! Instead of using y_{noise}

$$e = y_{noise} - g_{ideal}(\mathbf{x}) \quad (11)$$

We use y_{ideal}

$$e = y_{ideal} - g_{ideal}(\mathbf{x}) \quad (12)$$

We will see

How the geometry will solve the problem with using these labels.

Outline

1 Introduction

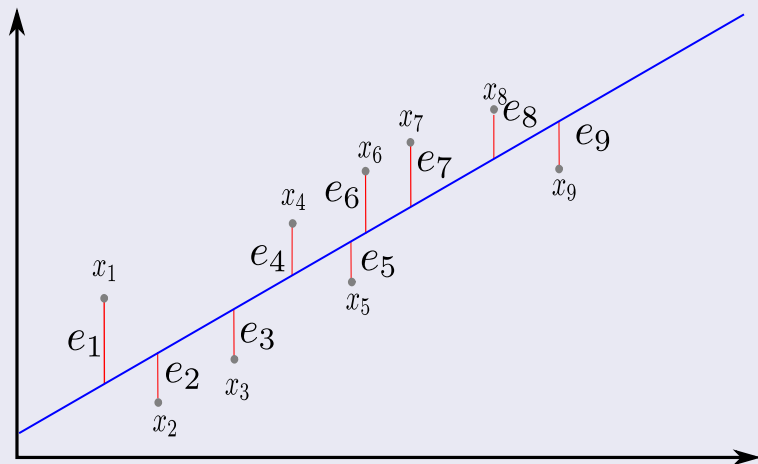
- Introduction
- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector

2 Developing a Solution

- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- **The Final Error Equation**
- Geometric Interpretation
- Issues with Least Squares!!!
 - Problem with Outliers
 - Problem with High Number of Dimensions

Here, we have multiple errors

What can we do?



Sum Over All the Errors

We can do the following

$$J(\mathbf{w}) = \sum_{i=1}^N e_i^2 = \sum_{i=1}^N (y_i - g_{ideal}(\mathbf{x}_i))^2 \quad (13)$$

Remark: This is known as the Least Squared Error cost function

Sum Over All the Errors

We can do the following

$$J(\mathbf{w}) = \sum_{i=1}^N e_i^2 = \sum_{i=1}^N (y_i - g_{ideal}(\mathbf{x}_i))^2 \quad (13)$$

Remark: This is known as the Least Squared Error cost function

Generalizing

- The dimensionality of each sample (data point) is d .

• You can extend each vector sample to be $\mathbf{x}^T = (1, \mathbf{x}')$.

Sum Over All the Errors

We can do the following

$$J(\mathbf{w}) = \sum_{i=1}^N e_i^2 = \sum_{i=1}^N (y_i - g_{ideal}(\mathbf{x}_i))^2 \quad (13)$$

Remark: This is known as the Least Squared Error cost function

Generalizing

- The dimensionality of each sample (data point) is d .
- You can extend each vector sample to be $\mathbf{x}^T = (\mathbf{1}, \mathbf{x}')$.

We can use a trick

The following function

$$g_{ideal}(\mathbf{x}) = \begin{pmatrix} 1 & x_1 & x_2 & \dots & x_d \end{pmatrix} \begin{pmatrix} w_0 \\ w_1 \\ w_2 \\ w_3 \\ \vdots \\ w_d \end{pmatrix} = \mathbf{x}^T \mathbf{w}$$

We can rewrite the error equation as

$$J(\mathbf{w}) = \sum_{i=1}^N (y_i - g_{ideal}(\mathbf{x}_i))^2 = \sum_{i=1}^N (y_i - \mathbf{x}_i^T \mathbf{w})^2 \quad (14)$$

We can use a trick

The following function

$$g_{ideal}(\mathbf{x}) = \begin{pmatrix} 1 & x_1 & x_2 & \dots & x_d \end{pmatrix} \begin{pmatrix} w_0 \\ w_1 \\ w_2 \\ w_3 \\ \vdots \\ w_d \end{pmatrix} = \mathbf{x}^T \mathbf{w}$$

We can rewrite the error equation as

$$J(\mathbf{w}) = \sum_{i=1}^N (y_i - g_{ideal}(\mathbf{x}_i))^2 = \sum_{i=1}^N (y_i - \mathbf{x}_i^T \mathbf{w})^2 \quad (14)$$

Furthermore

Then stacking all the possible estimations into the product Data Matrix and weight vector

$$\mathbf{X}\mathbf{w} = \begin{pmatrix} 1 & (\mathbf{x}_1)_1 & \cdots & (\mathbf{x}_1)_j & \cdots & (\mathbf{x}_1)_d \\ \vdots & & & \vdots & & \vdots \\ 1 & (\mathbf{x}_i)_1 & & (\mathbf{x}_i)_j & & (\mathbf{x}_i)_d \\ \vdots & & & \vdots & & \vdots \\ 1 & (\mathbf{x}_N)_1 & \cdots & (\mathbf{x}_N)_j & \cdots & (\mathbf{x}_N)_d \end{pmatrix} \begin{pmatrix} w_1 \\ w_2 \\ w_3 \\ \vdots \\ w_{d+1} \end{pmatrix}$$

Note about other representations

We could have $\mathbf{x}^T = (x_1, x_2, \dots, x_d, 1)$ thus

$$\mathbf{X} = \begin{pmatrix} (\mathbf{x}_1)_1 & \cdots & (\mathbf{x}_1)_j & \cdots & (\mathbf{x}_1)_d & 1 \\ & & \vdots & & \vdots & \vdots \\ (\mathbf{x}_i)_1 & & (\mathbf{x}_i)_j & & (\mathbf{x}_i)_d & 1 \\ & & \vdots & & \vdots & \vdots \\ (\mathbf{x}_N)_1 & \cdots & (\mathbf{x}_N)_j & \cdots & (\mathbf{x}_N)_d & 1 \end{pmatrix} \quad (15)$$

Then, we have the following trick with \mathbf{X}

With the Data Matrix

$$\mathbf{X}\mathbf{w} = \begin{pmatrix} \mathbf{x}_1^T \mathbf{w} \\ \mathbf{x}_2^T \mathbf{w} \\ \mathbf{x}_3^T \mathbf{w} \\ \vdots \\ \mathbf{x}_N^T \mathbf{w} \end{pmatrix} \quad (16)$$

Therefore

We have that

$$\begin{pmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_4 \end{pmatrix} - \begin{pmatrix} \mathbf{x}_1^T \mathbf{w} \\ \mathbf{x}_2^T \mathbf{w} \\ \mathbf{x}_3^T \mathbf{w} \\ \vdots \\ \mathbf{x}_N^T \mathbf{w} \end{pmatrix} = \begin{pmatrix} y_1 - \mathbf{x}_1^T \mathbf{w} \\ y_2 - \mathbf{x}_2^T \mathbf{w} \\ y_3 - \mathbf{x}_3^T \mathbf{w} \\ \vdots \\ y_4 - \mathbf{x}_N^T \mathbf{w} \end{pmatrix}$$

Then, we have the following equality:

$$\begin{pmatrix} y_1 - \mathbf{x}_1^T \mathbf{w} & y_2 - \mathbf{x}_2^T \mathbf{w} & y_3 - \mathbf{x}_3^T \mathbf{w} & \dots & y_4 - \mathbf{x}_N^T \mathbf{w} \end{pmatrix} \begin{pmatrix} y_1 - \mathbf{x}_1^T \mathbf{w} \\ y_2 - \mathbf{x}_2^T \mathbf{w} \\ y_3 - \mathbf{x}_3^T \mathbf{w} \\ \vdots \\ y_4 - \mathbf{x}_N^T \mathbf{w} \end{pmatrix} = \sum_{i=1}^N (y_i - \mathbf{x}_i^T \mathbf{w})^2$$

Therefore

We have that

$$\begin{pmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_4 \end{pmatrix} - \begin{pmatrix} \mathbf{x}_1^T \mathbf{w} \\ \mathbf{x}_2^T \mathbf{w} \\ \mathbf{x}_3^T \mathbf{w} \\ \vdots \\ \mathbf{x}_N^T \mathbf{w} \end{pmatrix} = \begin{pmatrix} y_1 - \mathbf{x}_1^T \mathbf{w} \\ y_2 - \mathbf{x}_2^T \mathbf{w} \\ y_3 - \mathbf{x}_3^T \mathbf{w} \\ \vdots \\ y_4 - \mathbf{x}_N^T \mathbf{w} \end{pmatrix}$$

Then, we have the following equality

$$\begin{pmatrix} y_1 - \mathbf{x}_1^T \mathbf{w} & y_2 - \mathbf{x}_2^T \mathbf{w} & y_3 - \mathbf{x}_3^T \mathbf{w} & \cdots & y_4 - \mathbf{x}_N^T \mathbf{w} \end{pmatrix} \begin{pmatrix} y_1 - \mathbf{x}_1^T \mathbf{w} \\ y_2 - \mathbf{x}_2^T \mathbf{w} \\ y_3 - \mathbf{x}_3^T \mathbf{w} \\ \vdots \\ y_4 - \mathbf{x}_N^T \mathbf{w} \end{pmatrix} = \sum_{i=1}^N (y_i - \mathbf{x}_i^T \mathbf{w})^2$$

Then, we have

The following equality

$$\sum_{i=1}^N \left(y_i - \mathbf{x}_i^T \mathbf{w} \right)^2 = (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) = \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 \quad (17)$$

The Final Discriminant Function

Very Simple!!!

$$g(\mathbf{x}) = \mathbf{x}^T \mathbf{w} = \mathbf{x}^T \left(\mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T \mathbf{y} \quad (18)$$

Pseudo-inverse of a Matrix

Definition

Suppose that $X \in \mathbb{R}^{m \times n}$ and $\text{rank}(X) = m$. We call the matrix

$$X^+ = (X^T X)^{-1} X^T$$

the pseudo inverse of X .

Pseudo-inverse of a Matrix

Definition

Suppose that $X \in \mathbb{R}^{m \times n}$ and $\text{rank}(X) = m$. We call the matrix

$$\mathbf{X}^+ = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$$

the pseudo inverse of X .

Pseudo-inverse of a Matrix

Definition

Suppose that $X \in \mathbb{R}^{m \times n}$ and $\text{rank}(X) = m$. We call the matrix

$$\mathbf{X}^+ = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$$

the pseudo inverse of X .

Pseudo-inverse of a Matrix

Definition

Suppose that $X \in \mathbb{R}^{m \times n}$ and $\text{rank}(X) = m$. We call the matrix

$$\mathbf{X}^+ = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$$

the pseudo inverse of X .

Outline

1 Introduction

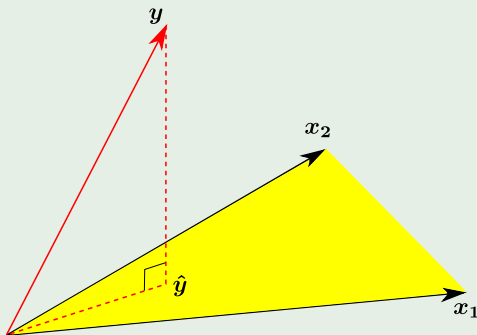
- Introduction
- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector

2 Developing a Solution

- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation
- **Geometric Interpretation**
- Issues with Least Squares!!!
 - Problem with Outliers
 - Problem with High Number of Dimensions

Geometrically

Given a y , you obtain a projected \hat{y} through the process $X^T y$



This Resolve Our Problem

With the Labels being chosen at the beginning

Question? Did you noticed the following?

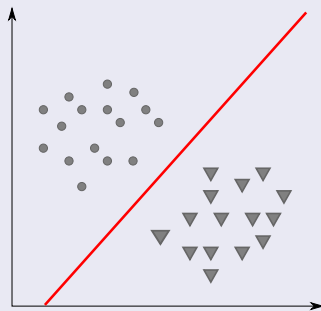
We assume a similar number of elements in both classes

This Resolve Our Problem

With the Labels being chosen at the beginning

Question? Did you noticed the following?

We assume a similar number of elements in both classes



Outline

1 Introduction

- Introduction
- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector

2 Developing a Solution

- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation
- Geometric Interpretation
- **Issues with Least Squares!!!**
 - Problem with Outliers
 - Problem with High Number of Dimensions

Outline

1 Introduction

- Introduction
- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector

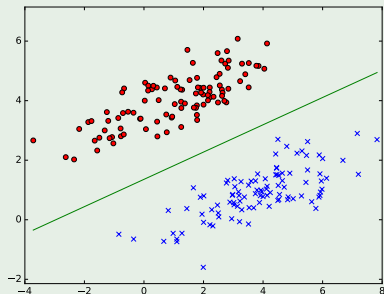
2 Developing a Solution

- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation
- Geometric Interpretation
- **Issues with Least Squares!!!**
 - **Problem with Outliers**
 - Problem with High Number of Dimensions

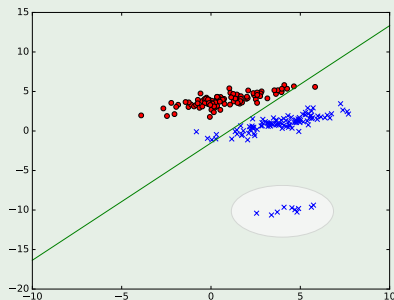
Issues with Least Squares

Problem with Outliers

No Outliers



Outliers



Outline

1 Introduction

- Introduction
- The Simplest Functions
- Splitting the Space
- Defining the Decision Surface
- Properties of the Hyperplane $w^T x + w_0$
- Augmenting the Vector

2 Developing a Solution

- Least Squared Error Procedure
 - The Geometry of a Two-Category Linearly-Separable Case
- The Error Idea
- The Final Error Equation
- Geometric Interpretation
- **Issues with Least Squares!!!**
 - Problem with Outliers
 - Problem with High Number of Dimensions

Problems with a High Number of Dimensions

In Many Modern Problems

- Many dimensions/features/predictors (possibly thousands).

Problems with a High Number of Dimensions

In Many Modern Problems

- Many dimensions/features/predictors (possibly thousands).

Only a few of these may be important

- It needs some form of feature selection.

• Possible some type of regularization.

Problems with a High Number of Dimensions

In Many Modern Problems

- Many dimensions/features/predictors (possibly thousands).

Only a few of these may be important

- It needs some form of feature selection.
- Possible some type of regularization.

Why?

- Least Square Error Regression treats all dimensions equally.
- Relevant dimensions might be averaged with irrelevant ones.

Problems with a High Number of Dimensions

In Many Modern Problems

- Many dimensions/features/predictors (possibly thousands).

Only a few of these may be important

- It needs some form of feature selection.
- Possible some type of regularization.

Why?

- Least Square Error Regression treats all dimensions equally.

• Relevant dimensions might be averaged with irrelevant ones.

Problems with a High Number of Dimensions

In Many Modern Problems

- Many dimensions/features/predictors (possibly thousands).

Only a few of these may be important

- It needs some form of feature selection.
- Possible some type of regularization.

Why?

- Least Square Error Regression treats all dimensions equally.
- Relevant dimensions might be averaged with irrelevant ones.